



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Fundamentals and Recent Developments in Approximate Bayesian Computation

Citation for published version:

Lintusaari, J, Gutmann, M, Dutta, R, Kaski, S & Corander, J 2017, 'Fundamentals and Recent Developments in Approximate Bayesian Computation', *Systematic biology*, vol. 66, no. 1, pp. e66-e82.
<https://doi.org/10.1093/sysbio/syw077>

Digital Object Identifier (DOI):

[10.1093/sysbio/syw077](https://doi.org/10.1093/sysbio/syw077)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Systematic biology

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Fundamentals and Recent Developments in ABC

Fundamentals and Recent Developments in Approximate Bayesian Computation

JARNO LINTUSAARI^{1*}, MICHAEL U. GUTMANN^{1,2*}, RITABRATA DUTTA¹,

SAMUEL KASKI¹, AND JUKKA CORANDER^{2,3}

¹*Helsinki Institute for Information Technology HIIT,*

Department of Computer Science, Aalto University, Espoo, 00076, Finland

²*Helsinki Institute for Information Technology HIIT,*

Department of Mathematics and Statistics, University of Helsinki, Helsinki, 00014, Finland

³*Department of Biostatistics, University of Oslo, 0317 Oslo, Norway*

*Contributed equally

Corresponding author: Jarno Lintusaari, Department of Computer Science, Aalto University, 00076 Espoo, Finland, E-mail: jarno.lintusaari@aalto.fi

Abstract.— Bayesian inference plays an important role in phylogenetics, evolutionary biology and in many other branches of science. It provides a principled framework for dealing with uncertainty and quantifying how it changes in the light of new evidence. For many complex models and inference problems, however, only approximate quantitative answers are obtainable. Approximate Bayesian computation (ABC) refers to a family of algorithms for approximate inference that make a minimal set of assumptions by only requiring that sampling from a model is possible. We explain here the fundamentals of approximate Bayesian computation, review the classical algorithms, and highlight recent developments.

(Keywords: ABC, approximate Bayesian computation, Bayesian inference, likelihood-free inference, phylogenetics, simulator-based models, stochastic simulation models, tree-based models)

INTRODUCTION

Many recent models in biology describe nature to a high degree of accuracy but are not amenable to analytical treatment. The models can, however, be simulated on computers and we can thereby replicate many complex phenomena such as the evolution of genomes (Marttinen et al. 2015), the dynamics of gene regulation (Toni et al. 2009), or the demographic spread of a species (Currat and Excoffier 2004; Fagundes et al. 2007; Itan et al. 2009; Excoffier et al. 2013). Such simulator-based models are often stochastic and have multiple parameters. While it is usually relatively easy to generate data from the models for any configuration of the parameters, the real interest is often focused on the inverse problem: the identification of parameter configurations that would plausibly lead to data that are sufficiently similar to the observed data. Solving such a nonlinear inverse problem is generally a very difficult task.

Bayesian inference provides a principled framework for solving the aforementioned inverse problem. A prior probability distribution on the model parameters is used to describe the initial beliefs about what values of the parameters could be plausible. The prior beliefs are updated in light of the observed data by means of the likelihood function. Computing the likelihood function, however, is mostly impossible for simulator-based models due to the unobservable (latent) random quantities that are present in the model. In some cases, Monte Carlo methods offer a way to handle the latent variables such that an approximate likelihood is obtained, but these methods have their limitations, and for large and complex models, they are “too inefficient by far” (Green et al. 2015, page 848). To deal with models where likelihood calculations fail, other techniques have been developed that are collectively referred to as likelihood-free inference or approximate Bayesian computation (ABC).

In a nutshell, ABC algorithms sample from the posterior distribution of the

parameters by finding values that yield simulated data sufficiently resembling the observed data. ABC is widely used in systematics. For instance, Hickerson et al. (2006) used ABC to test for simultaneous divergence between members of species pairs. Fan and Kubatko (2011) estimated the topology and speciation times of a species tree under the coalescent model using ABC. Their method does not require sequence data, only gene tree topology information, and was found to perform favorably in terms of both accuracy and computation time. Slater et al. (2012) used ABC to simultaneously infer rates of diversification and trait evolution from incompletely sampled phylogenies and trait data. They found their ABC approach to be comparable to likelihood-based methods that use complete datasets. In addition, the ABC approach can handle extremely sparsely sampled phylogenies and trees containing very large numbers of species. Ratmann et al. (2012) used ABC to fit two different mechanistic phylodynamic models for interpandemic influenza A(H3N2) using both surveillance data and sequence data simultaneously. The simultaneous consideration of these two types of data allowed them to drastically constrain the parameter space and expose model deficiencies using the ABC framework. Very recently Baudet et al. (2015) used ABC to reconstruct the coevolutionary history of host-parasite systems. The ABC-based method was shown to handle large trees beyond the scope of other existing methods.

While widely applicable, ABC comes with its own set of difficulties, that are of both computational and statistical nature. The two main intrinsic difficulties are how to efficiently find plausible parameter values, and how to define what is similar to the observed data and what is not. All ABC algorithms have to deal with these two issues in some manner, and the different algorithms discussed here essentially differ in how they tackle the two problems.

The remainder of this article is structured as follows. We next discuss important properties of simulator-based models and point out difficulties when performing statistical

inference with them. The discussion leads to the basic rejection ABC algorithm which is presented in the subsequent section. This is followed by a presentation of popular ABC algorithms that have been developed to increase the computational efficiency. We then consider several recent advances that aim to improve ABC both computationally and statistically. The final section provides conclusions and a discussion about likelihood-free inference methods related to ABC.

SIMULATOR-BASED MODELS

Definition

Simulator-based models are functions M that map the model parameters θ and some random variables V to data y . The functions M are generally implemented as computer programs where the parameter values are provided as input and where the random variables are drawn sequentially by making calls to a random number generator. The parameters θ govern the properties of interest of the generated data while the random variables V represent the stochastic variation inherent to the simulated process.

The mapping M may be as complex as needed, and this generality of simulator-based models allows researchers to implement hypotheses about how the data were generated without having to make excessive compromises motivated by mathematical simplicity, or other reasons not related to the scientific question being investigated.

Due to the presence of the random variables V , the outputs of the simulator fluctuate randomly even when using exactly the same values of the model parameters θ . This means that we can consider the simulator to define a random variable Y_θ whose distribution is implicitly determined by the distribution of V and the mapping M acting on V for a given θ (for this reason, simulator-based models are sometimes called implicit

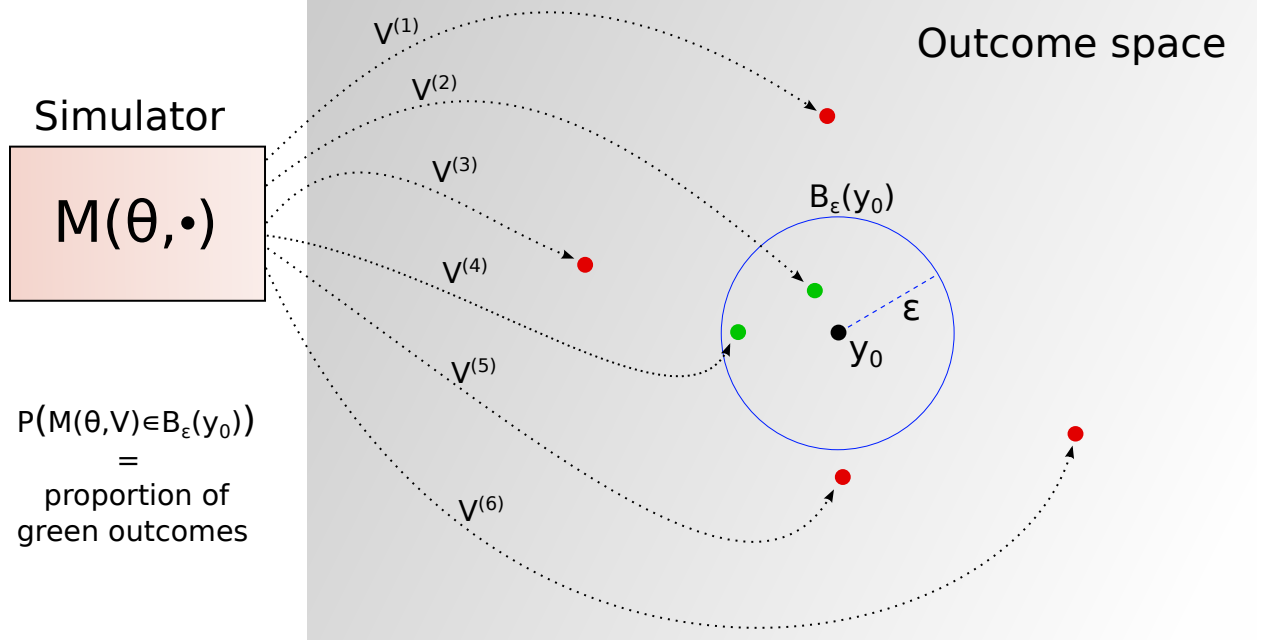


Figure 1: Illustration of the stochastic simulator M run multiple times with a fixed value of θ . The black dot y_0 is the observed data and the arrows point to different simulated data sets. Two outcomes, marked in green, are less than ϵ away from y_0 . The proportion of such outcomes provides an approximation of the likelihood of θ for observed data y_0 .

models, Diggle and Gratton 1984). Using the properties of transformation of random variables, it is possible to formally write down the distribution of Y_θ . For instance, for a fixed value of θ , the probability that Y_θ takes values in an ϵ neighborhood $B_\epsilon(y_0)$ around the observed data y_0 is equal to the probability to draw values of V that are mapped to that neighborhood (Figure 1),

$$\Pr(Y_\theta \in B_\epsilon(y_0)) = \Pr(M(\theta, V) \in B_\epsilon(y_0)). \quad (1)$$

Computing the probability analytically is impossible for complex models. But it is possible to test empirically whether a particular outcome y_θ of the simulation ends up in the neighborhood of y_0 or not (see Figure 1). We will see that this property of simulator-based models plays a key role in performing inference about their parameters.

Example

As an example of a simulator-based model, we here present the simple yet analytically intractable model by Tanaka et al. (2006) for the spread of tuberculosis. We will use the model throughout the paper for illustrating different concepts and methods.

The model begins with one infectious host and stops when a fixed number of infectious hosts m is exceeded (Figure 2). In the simulation, it is assumed that each infectious host randomly infects other individuals from an unlimited supply of hosts with the rate α , each time transmitting a specific strain of the communicable pathogen, characterized by its haplotype. It is thus effectively assumed that a strong transmission bottleneck occurs, such that only a single strain is passed forward in each transmission event, despite the eventual genetic variation persisting in the within-host pathogen population. Further, each infected host is considered to be infectious immediately. The model states that a host stops being infectious, i.e. recovers or dies, randomly with the rate δ , and the pathogen of the host mutates randomly within the host at the rate τ , thereby generating a novel haplotype under a single-locus infinite alleles model. The parameters of the model are thus $\theta = (\alpha, \delta, \tau)$. The output of the simulator is a vector of cluster sizes in the simulated population of infected hosts, where clusters are the groups of hosts infected by the same haplotype of the pathogen. After the simulation, a random sample of size $n < m$ is taken from the population yielding the vector of cluster sizes y_θ present in the sample. For example, $y_\theta = (6, 3, 2, 2, 1, 1, 1, 1, 1, 1)$ corresponds to a sample of size 20 containing one cluster with 6 infected hosts, one cluster with three hosts, two clusters with two hosts each, as well as 7 singleton clusters. Note that this model of pathogen spread is atypical in the sense that the observation times of the infections are all left implicit in the sampling process, in contrast to the standard likelihood formulation used for infectious disease epidemiological models (Anderson and May 1992).

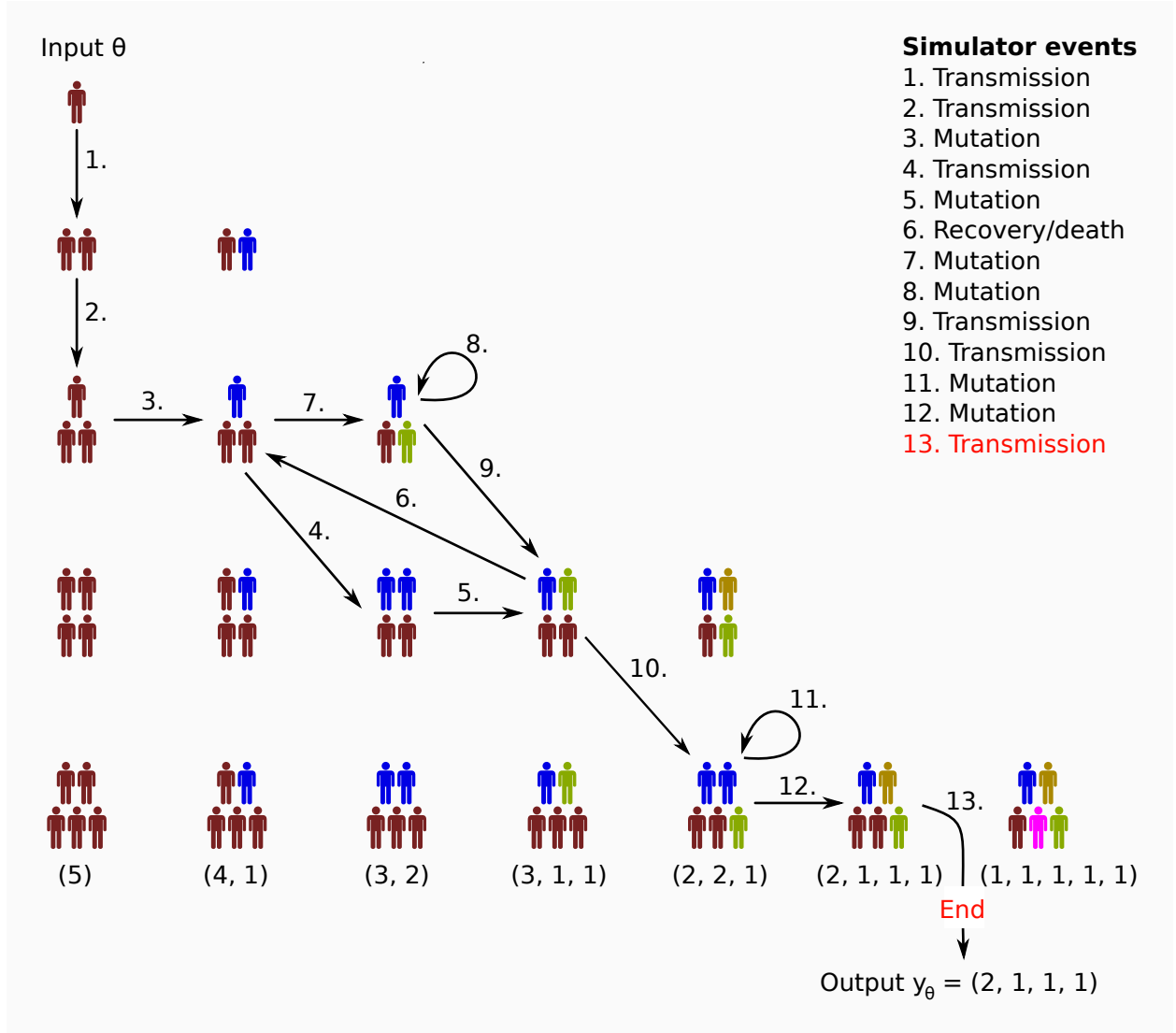


Figure 2: An example of a transmission process simulated under a parameter configuration θ without sub-sampling of the simulated infectious population. Arrows indicate the sequence of random events taking place in the simulation and different colors represent different haplotypes of the pathogen. The simulation starts with one infectious host who transmits the pathogen to another host. After one more transmission event, the pathogen undergoes a mutation within one of the three hosts infected so far (event three). As the sixth event in the simulation, one of the haplotypes is removed from the population due to the recovery/death of the corresponding host. The simulation stops when the infectious population size exceeds $m = 5$ and the simulator outputs the generated y_θ . The nodes not connected by arrows show all the other possible configurations of the infectious population, but which were not visited in this example run of the simulator. The bottom row lists the possible outputs of the simulator (cluster size vectors) under their corresponding population configuration.

Difficulties in performing statistical inference

Values of the parameters θ that are plausible in the light of the observations y_0 can be determined via statistical inference either by finding values that maximize the probability in Equation (1) for some sufficiently small ϵ , or by determining their posterior distribution. In more detail, in maximum likelihood estimation, the parameters are determined by maximizing the likelihood function $L(\theta)$,

$$L(\theta) = \lim_{\epsilon \rightarrow 0} c_\epsilon \Pr(Y_\theta \in B_\epsilon(y_0)), \quad (2)$$

where c_ϵ is a proportionality factor that may depend on ϵ , which is needed when $\Pr(Y_\theta \in B_\epsilon(y_0))$ shrinks to zero as ϵ approaches zero. If the output of the simulator can only take a countable number of values, Y_θ is called a discrete random variable and the definition of the likelihood simplifies to $L(\theta) = \Pr(Y_\theta = y_0)$, which equals the probability of simulating data equal to the observed data. In Bayesian inference, the essential characterization of the uncertainty about the model parameters is defined by their conditional distribution given the data, i.e. the posterior distribution $p(\theta|y_0)$,

$$p(\theta|y_0) \propto L(\theta)p(\theta), \quad (3)$$

where $p(\theta)$ is the prior distribution of the parameters.

For complex models neither the probability in Equation (1), nor the likelihood function $L(\theta)$ are available analytically in closed form as a function of θ , which is the reason why statistical inference is difficult for simulator-based models.

For the model of tuberculosis transmission presented in the previous section, computing the likelihood function becomes intractable if the infectious population size m is large, or if the death rate $\delta > 0$ (Stadler 2011). This is because for large m , the state space

Algorithm 1 Rejection sampling algorithm for simulator-based models. The algorithm produces N independent samples $\theta^{(i)}$ from the posterior distribution $p(\theta|y_0)$

```

1: for  $i = 1$  to  $N$  do
2:   repeat
3:     Generate  $\theta$  from the prior  $p(\cdot)$ 
4:     Generate  $y_\theta$  from the simulator
5:   until  $y_\theta = y_0$ 
6:    $\theta^{(i)} \leftarrow \theta$ 
7: end for

```

of the process, i.e. the number of different cluster vectors, grows very quickly. This makes exact numerical calculation of the likelihood infeasible because in essence, every possible path to the outcome should be accounted for (Figure 2). Moreover, if the death rate δ is nonzero, the process is allowed to return to previous states which further complicates the computations. Finally, the assumption that not all infectious hosts are observed contributes additionally to the intractability of the likelihood. Stadler (2011) approached the problem using transmission trees (Figure 3). The likelihood function stays, however, intractable because of the vast number of different trees that all yield the same observed data and thus need to be considered when evaluating the likelihood of a parameter value.

Inference via rejection sampling

We present here an algorithm for exact posterior inference that is applicable when Y_θ can only take countably many values, that is, if Y_θ is a discrete random variable. As shown above, in this case $L(\theta) = \Pr(Y_\theta = y_0)$. The presented algorithm forms the basis of the algorithms for approximate Bayesian computation discussed in the later sections.

In general, samples from the prior distribution $p(\theta)$ of the parameters can be converted into samples from the posterior $p(\theta|y_0)$ by retaining each sampled value with a probability proportional to $L(\theta)$. This can be done sequentially by first sampling a parameter value from the prior, $\theta \sim p(\theta)$, and then accepting the obtained value with the

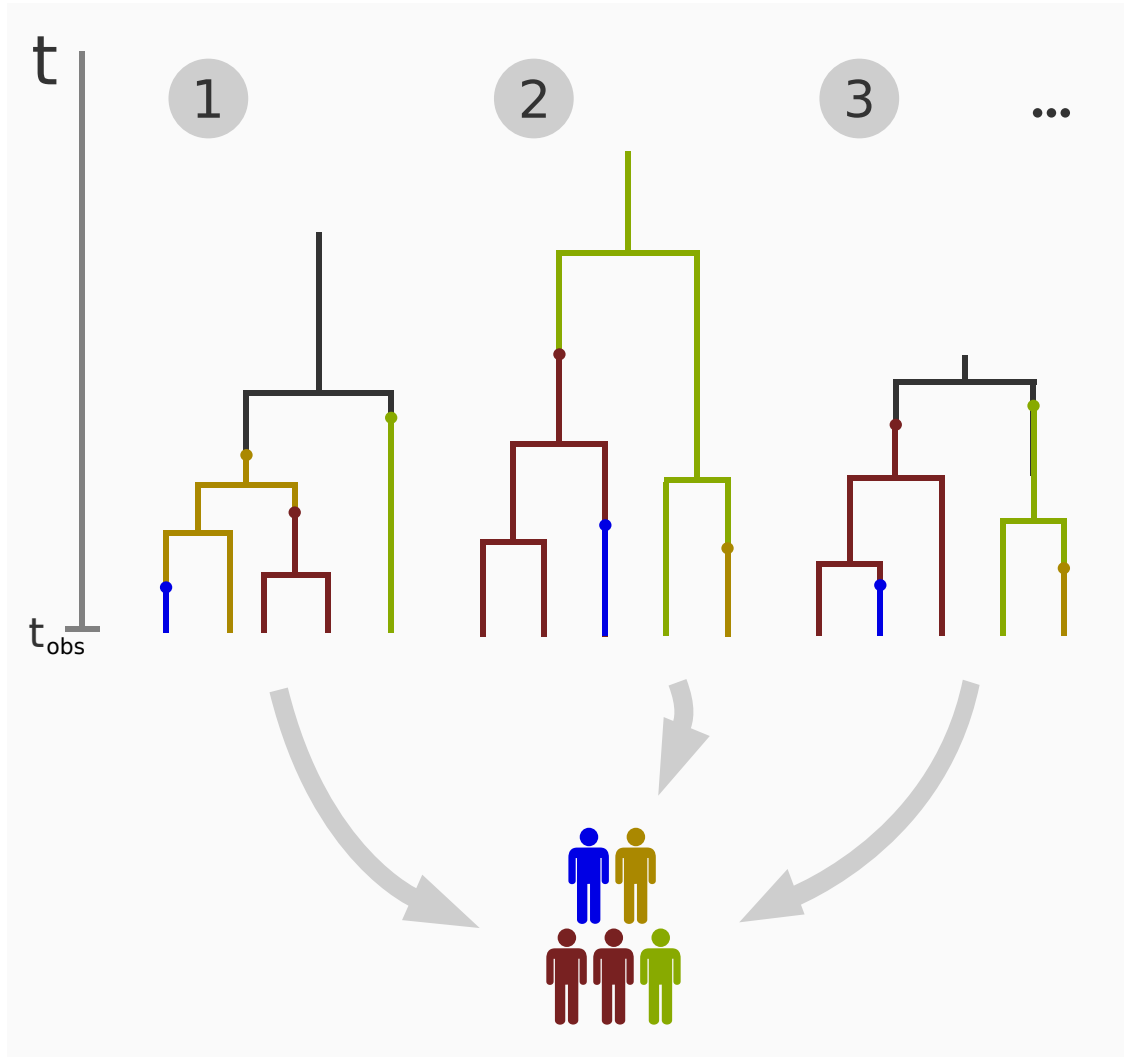


Figure 3: The transmission process in Figure 2 can also be described with transmission trees (Stadler 2011) paired with mutations. The trees are characterized by their structure, the length of their edges, and the mutations on the edges (marked with small circles that change the color of the edge, where colors represent the different haplotypes of the pathogen). The figure shows three examples of different trees that yield the same observed data at the observation time t_{obs} . Calculating the likelihood of a parameter value requires summing over all possible trees yielding the observed data, which is computationally impossible when the sample size is large.

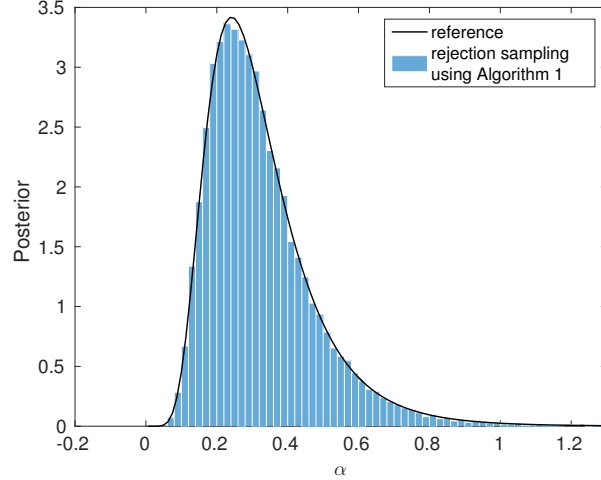


Figure 4: Exact inference for a simulator-based model of tuberculosis transmission. A very simple setting was chosen where the exact posterior can be numerically computed (black line), and where Algorithm 1 is applicable (blue bars).

probability $L(\theta)/(\max_{\theta} L(\theta))$. This procedure corresponds to rejection sampling (see, for example Robert and Casella 2004, Chapter 2). Now, with the likelihood $L(\theta)$ being equal to the probability that $Y_{\theta} = y_0$, the latter step can be implemented for simulator-based models even when $L(\theta)$ is not available analytically: we run the simulator and check whether the generated data equal the observed data. This gives the rejection algorithm for simulator-based models summarized as Algorithm 1. Rubin (1984) used it to provide intuition about how Bayesian inference about parameters works in general.

To obtain another interpretation of Algorithm 1, recall that for discrete random variables the posterior distribution $p(\theta|y_0)$ is, by definition, equal to the joint distribution of θ and Y_{θ} , normalized by the probability that $Y_{\theta} = y_0$. That is, the posterior is obtained by conditioning on the event $Y_{\theta} = y_0$. We can thus understand the test for equality of y_{θ} and y_0 on line 5 of the algorithm as an implementation of the conditioning operation.

To illustrate Algorithm 1, we generated a synthetic data set y_0 from the tuberculosis transmission model by running the simulator with the parameter values $\alpha = 0.2$, $\delta = 0$, $\tau = 0.198$, and setting the population size to $m = 20$. We further assumed that the whole

population is observed, which yielded the observed data $y_0 = (6, 3, 2, 2, 1, 1, 1, 1, 1, 1, 1)$. The assumptions about the size of the population, and that the whole population was observed, are unrealistic but they enable a comparison to the exact posterior distribution, which in this setting can be numerically computed using Theorem 1 of Stadler (2011). In this case, the histogram of samples obtained with Algorithm 1 matches the posterior distribution very accurately (Figure 4). To obtain this result, we assumed that both of the parameters δ and τ were known, and assigned a uniform prior distribution in the interval $(0.005, 2)$ for the sole unknown parameter, the transmission rate α . A total of 20 million data sets y_θ were simulated, out of which 40 000 matched y_0 (acceptance rate of 0.2%).

FUNDAMENTALS OF APPROXIMATE BAYESIAN COMPUTATION

The rejection ABC algorithm

While Algorithm 1 produces independent samples from the posterior, the probability that the simulated data equal the observed data is often negligibly small, which renders the algorithm impractical as virtually no simulated realizations of θ will be accepted. The same problem holds true if the generated data can take uncountably many values, i.e. when Y_θ is a continuous random variable.

To make inference feasible, the acceptance criterion $y_\theta = y_0$ in Algorithm 1 can be relaxed to

$$d(y_\theta, y_0) \leq \epsilon, \tag{4}$$

where $\epsilon > 0$ and $d(y_\theta, y_0) \geq 0$ is a “distance” function that measures the discrepancy between the two data sets, as considered relevant for the inference. With this modification, Algorithm 1 becomes the rejection ABC algorithm summarized as Algorithm 2. The first implementation of this algorithm appeared in the work by Pritchard et al. (1999).

Algorithm 2 does not produce samples from the posterior $p(\theta|y_0)$ in Equation (3) but samples from an approximation $p_{d,\epsilon}(\theta|y_0)$,

$$p_{d,\epsilon}(\theta|y_0) \propto \Pr(d(Y_\theta, y_0) \leq \epsilon)p(\theta), \quad (5)$$

which is the posterior distribution of θ conditional on the event $d(Y_\theta, y_0) \leq \epsilon$. Equation (5) is obtained by approximating the intractable likelihood function $L(\theta)$ in Equation (2) with $L_{d,\epsilon}(\theta)$,

$$L_{d,\epsilon}(\theta) \propto \Pr(d(Y_\theta, y_0) \leq \epsilon). \quad (6)$$

The approximation is two-fold. First, the distance function d is generally not a metric in the mathematical sense, namely $d(y_\theta, y_0) = 0$ even if $y_\theta \neq y_0$. This may happen, for example, when d is defined through summary statistics that remove information from the data (see below). Second, ϵ is chosen large enough so that enough samples will be accepted. Intuitively, the likelihood of a parameter value is approximated by the probability that running the simulator with said parameter value produces data within ϵ distance of y_0 (see Figure 1).

The distance d is typically computed by first reducing the data to suitable summary statistics $t = T(y)$ and then computing the distance d_T between them, so that $d(y_\theta, y_0) = d_T(t, t_0)$, where d_T is often the Euclidean or some other metric for the summary statistics. When combining different summary statistics, they are usually re-scaled so that they contribute equally to the distance (as, for example, done by Pritchard et al. 1999).

In addition to the accuracy of the approximation $p_{d,\epsilon}(\theta|y_0)$, the distance d and the threshold ϵ also influence the computing time required to obtain samples. For instance, if $\epsilon = 0$ and the distance d is such that $d(y, y_0) = 0$ if and only if $y = y_0$, then Algorithm 2 becomes Algorithm 1 and $p_{d,\epsilon}(\theta|y_0)$ becomes $p(\theta|y_0)$ but the computing time to obtain samples from $p_{d,\epsilon}(\theta|y_0)$ would typically be impractically large. Hence, on a very

Algorithm 2 Rejection ABC algorithm producing N independent samples from the approximate posterior distribution $p_{d,\epsilon}(\theta|y_0)$

```

1: for  $i = 1$  to  $N$  do
2:   repeat
3:     Generate  $\theta$  from the prior  $p(\cdot)$ 
4:     Generate  $y_\theta$  from the simulator
5:   until  $d(y_\theta, y_0) \leq \epsilon$ 
6:    $\theta^{(i)} \leftarrow \theta$ 
7: end for

```

fundamental level, there is a trade-off between statistical and computational efficiency in approximate Bayesian computation (see e.g. Beaumont et al. 2002, p. 2027).

We next illustrate Algorithm 2 and the mentioned trade-off using the previous example about tuberculosis transmission. Two distances d_1 and d_2 are considered,

$$d_1(y_\theta, y_0) = |T_1(y_\theta) - T_1(y_0)|, \quad d_2(y_\theta, y_0) = |T_2(y_\theta) - T_2(y_0)|, \quad (7)$$

where T_1 is the number of clusters contained in the data divided by the sample size n and T_2 is a genetic diversity measure defined as $T_2(y) = 1 - \sum_i (n_i/n)^2$, where n_i is the size of the i th cluster. For $y_0 = (6, 3, 2, 2, 1, 1, 1, 1, 1, 1)$, we have $T_1(y_0) = 11/20 = 0.55$ and $T_2(y_0) = 0.85$. For both d_1 and d_2 , the absolute difference between the summary statistics is used as the metric d_T .

For this example, using the summary statistic T_1 instead of the full data does not lead to a visible deterioration of the inferred posterior when $\epsilon = 0$ (Figure 5a). For summary statistic T_2 , however, there is a clear difference as the posterior mode and mean are shifted to larger values of α , and the posterior variance is larger too (Figure 5b). In both cases, increasing ϵ , that is, accepting more parameters, leads to an approximate posterior distribution that is less concentrated than the true posterior.

Algorithm 2 with summary statistic T_1 produces results comparable to Algorithm 1 but from the computational efficiency point of view the number of simulations required to

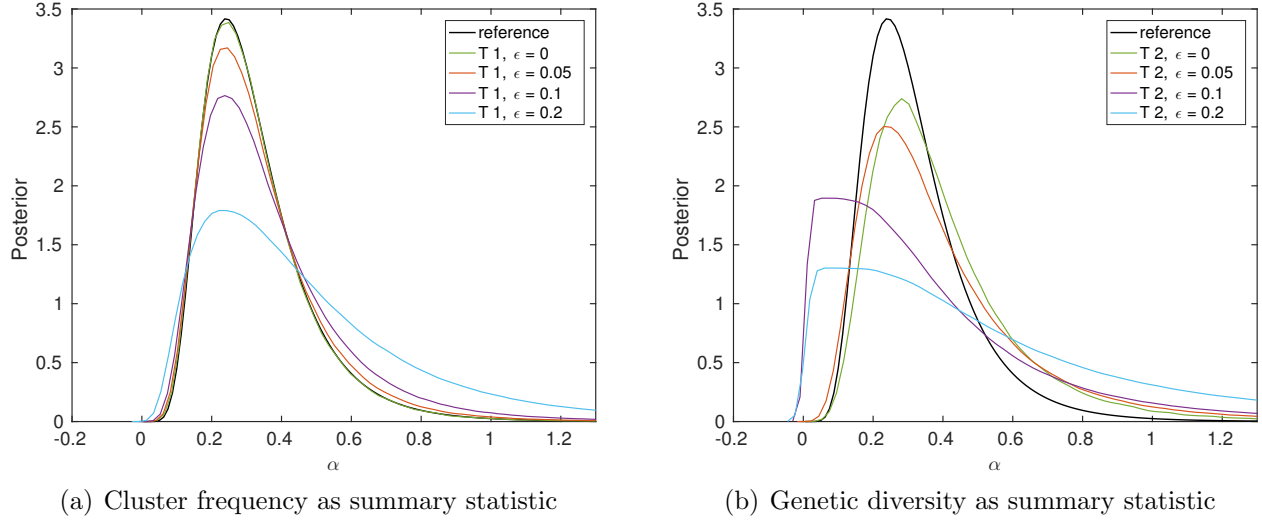


Figure 5: Inference results for the transmission rate α of tuberculosis. The plots show the posterior distributions obtained with Algorithm 2 and 20 million simulated data sets (proposals).

obtain the approximate posterior differs between the two algorithms. It can be seen that for a computational budget of 100 000 simulations, the posterior obtained by Algorithm 1 differs substantially from the exact posterior, while the posterior from Algorithm 2 with T_1 is still matching it well (Figure 6a). The relatively poor result with Algorithm 1 is due to its low acceptance rate (here 0.2%). While the accepted samples do follow the exact posterior $p(\theta|y_0)$, the algorithm did not manage to produce enough accepted realizations within the computational budget available, which implies that the Monte Carlo error of the posterior approximation remains non-negligible.

Plotting the number of data sets simulated versus the accuracy of the inferred posterior distribution allows us to further study the trade-off between statistical and computational efficiency between the different algorithms (Figure 6b). The accuracy is measured by the Kullback-Leibler (KL) divergence (Kullback and Leibler 1951) between the exact and the inferred posterior. Algorithm 2 with summary statistic T_1 features the best trade-off while using summary statistic T_2 performs the worst. The curve of the latter

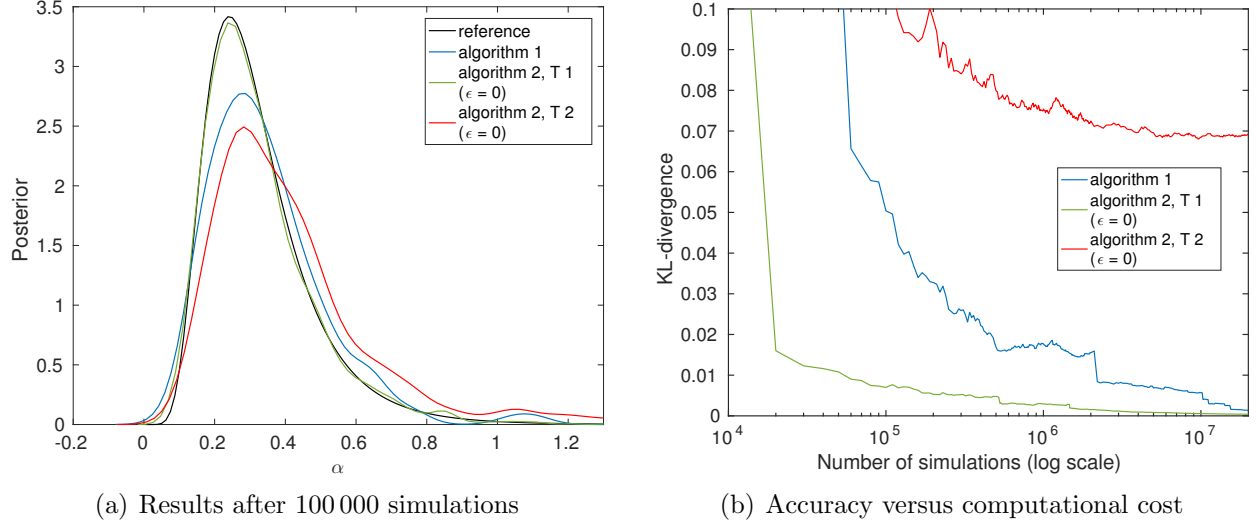


Figure 6: Comparison of the efficiency of Algorithms 1 and 2. Smaller Kullback-Leibler divergence means more accurate inference of the posterior distribution. Note that the stopping criterion of the algorithm has here been changed to be the total number of runs of the simulator instead of the number of accepted samples.

one flattens out after approximately one million simulations, showing the approximation error introduced by using the summary statistic T_2 . For Algorithm 1, nonzero values of the KL divergence are due to the Monte Carlo error only and it will approach the true posterior as the number of simulations grows. When using summary statistics, nonzero values of the KL divergence are due to both the Monte Carlo error and the use of the summary statistics. In this particular example, the error caused by the summary statistic T_1 is however negligible.

Choice of the summary statistics

If the distance d is computed by projecting the data to summary statistics followed by their comparison using a metric in the summary statistics space (e.g. the Euclidean distance), the quality of the inference hinges on the summary statistics chosen (Figures 5 and 6).

For consistent performance of ABC algorithms, the summary statistics should be

sufficient for the parameters, but this is often not the case. Additionally, with the increase in the number of summary statistics used, more simulations tend to be rejected so that an increasing number of simulation runs is needed to obtain a satisfactory number of accepted parameter values, making the algorithm computationally extremely inefficient. This is known as the curse of dimensionality for ABC (see also the discussion in the review paper by Beaumont 2010).

One of the main remedies to the above issue is to efficiently choose informative summary statistics. Importantly, the summary statistics that are informative for the parameters in a neighborhood of the true parameter value, and the summary statistics most informative globally, are significantly different (Nunes and Balding 2010). General intuition suggests that the set of summary statistics that are locally sufficient would be a subset of the globally sufficient ones. Therefore, a good strategy seems to first find a locality containing the true parameter with high enough probability and then choose informative statistics depending on that locality. However, this can be difficult in practice because rather different parameter values can produce summary statistics that are contained in the same locality.

In line with the above, Nunes and Balding (2010), Fearnhead and Prangle (2012), and Aeschbacher et al. (2012) first defined “locality” through a pilot ABC run and then chose the statistics in that locality. Four methods for choosing the statistics are generally used: a) a sequential scheme based on the principle of approximate sufficiency (Joyce and Marjoram 2008); b) selection of a subset of the summary statistics maximizing pre-specified criteria such as the Akaike information criterion (used by Blum et al. 2013) or the entropy of a distribution (used by Nunes and Balding 2010); c) partial least square regression which finds linear combinations of the original summary statistics that are maximally decorrelated and at the same time highly correlated with the parameters (Wegmann et al. 2009); d) assuming a statistical model between parameters and

transformed statistics of simulated data, summary statistics are chosen by minimizing a loss function (Fearnhead and Prangle 2012; Aeschbacher et al. 2012). For comparison of the above methods in simulated and practical examples, we refer readers to the work by Blum and François (2010), Aeschbacher et al. (2012), and Blum et al. (2013).

Choice of the threshold

Having the distance function d specified, possibly using summary statistics, the remaining factor in the approximation of the posterior in Equation (5) is the specification of the threshold ϵ .

Larger values of ϵ result in biased approximations $p_{d,\epsilon}(\theta|y_0)$ (see e.g. Figure 5). The gain is a faster algorithm, meaning a reduced Monte Carlo error as one is able to produce more samples per unit of time. Therefore, when specifying the threshold the goal is to find a good balance between the bias and the Monte Carlo error. We illustrate this using Algorithm 2 with the full data without reduction to summary statistics (in other words, $T(y) = y$). In this case, Algorithm 2 with $\epsilon = 0$ is identical to Algorithm 1. Here the choice $\epsilon = 3$ results in a better posterior compared to $\epsilon = 0$ when using a maximal number of 100 000 simulations (Figure 7a). This means that the gain from reduced Monte Carlo error is greater than the loss incurred by the bias. But this is no longer true for $\epsilon = 5$ where the bias dominates. Eventually, the exact method will converge to the true posterior, while the other two continue to suffer from the bias caused by the larger threshold (Figure 7b). However, with smaller computational budgets (less than 2 million simulations in our example), more accurate results are obtained with the nonzero threshold $\epsilon = 3$.

The choice of the threshold is typically made by experimenting with a precomputed pool of simulation-parameter pairs (y_θ, θ) . Rather than setting the threshold value by hand, it is often determined by accepting some small proportion of the simulations (e.g. 1%, see Beaumont et al. 2002). The choice between different options can be made more

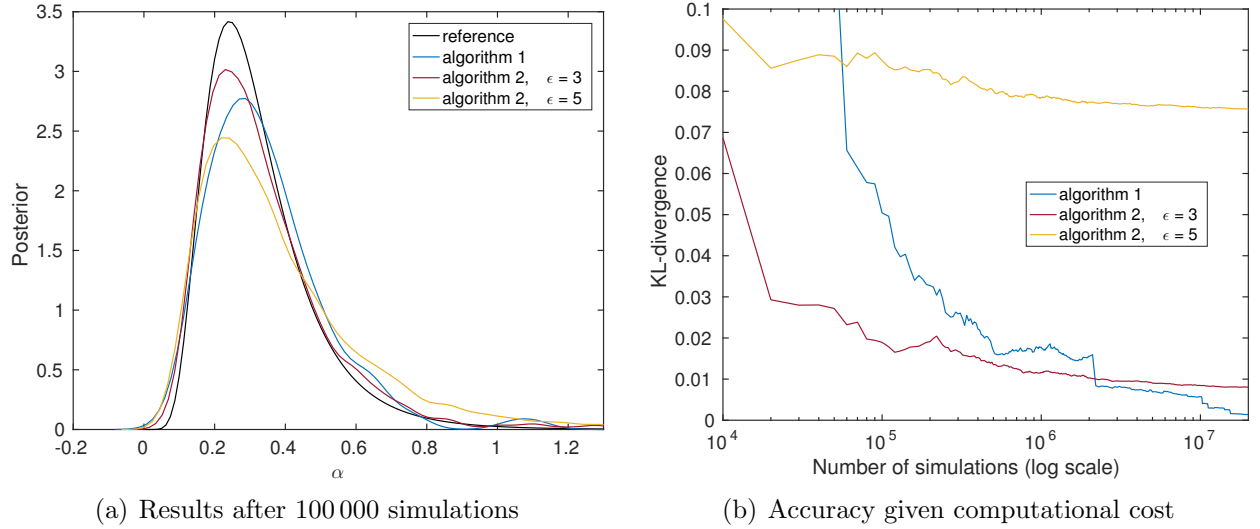


Figure 7: Comparison of the trade-off between Monte Carlo error and bias. Algorithm 1 is equivalent here to Algorithm 2 with $\epsilon = 0$. Smaller Kullback-Leibler divergences mean more accurate inference of the posterior distribution.

rigorous by using some of the simulated data sets in the role of the observed data and solving the inference problem for them using the remaining data sets. As the data-generating parameters are known for the simulated observations, different criteria, such as the mean squared error (MSE) between the mean of the approximation and the generating parameters can be used to make the choice (see e.g. Faisal et al. (2013), and the section on validation of ABC). This also allows one to assess the reliability of the inference procedure. Prangle et al. (2014) discuss the use of the coverage property (Wegmann et al. 2009) as the criterion to choose the threshold value ϵ . Intuitively, the coverage property tests if the parameter values θ^* used to artificially generate a data set y_0^* are covered by the credible intervals constructed from the ABC output for y_0^* at correct rates (i.e. $\alpha\%$ credible intervals should contain the true parameter in $\alpha\%$ of the tests).

If one plans to increase the computational budget after initial experiments, some of the theoretical convergence results can be used to adjust the threshold value. Barber et al. (2015) provide convergence results for an optimal ϵ sequence with respect to the mean

squared error of a posterior expectation (e.g. the posterior mean). The theoretically optimal sequence for the threshold ϵ is achieved by making it proportional to $N^{-1/4}$ as $N \rightarrow \infty$, where N is the number of accepted samples. If the constant in this relation is estimated in a pilot run, one can compute the new theoretically optimal threshold based on the planned increase in the computational budget. Blum (2010) derives corresponding results using an approach based on conditional density estimation, finding that ϵ should optimally be proportional to $N_s^{-1/(d+5)}$ as $N_s \rightarrow \infty$, where d is the dimension of the parameter space and N_s the total number of simulations performed (see also Fearnhead and Prangle (2012), Silk et al. (2013) and Biau et al. (2015) for similar results).

BEYOND SIMPLE REJECTION SAMPLING

The basic rejection ABC algorithm is essentially a trial and error scheme where the trial (proposal) values are sampled from the prior. We now review three popular algorithms that seek to improve upon the basic rejection approach. The first two aim at constructing proposal distributions that are closer to the posterior, whereas the third is a correction method that aims at adjusting samples obtained by ABC algorithms so that they are closer to the posterior.

Markov chain Monte Carlo ABC

The Markov chain Monte Carlo (MCMC) ABC algorithm is based on the Metropolis-Hastings MCMC algorithm which is often used in Bayesian statistics (Robert and Casella 2004, Chapter 7). In order to leverage this algorithm, we write $p_{d,\epsilon}(\theta|y_0)$ in Equation (5) as the marginal distribution of $p_{d,\epsilon}(\theta, y|y_0)$,

$$p_{d,\epsilon}(\theta, y|y_0) \propto p(\theta)p(y|\theta) \mathbb{1}[d(y, y_0) \leq \epsilon], \quad (8)$$

where $p(y|\theta)$ denotes the probability density (mass) function of Y_θ , and $\mathbb{1}[d(y, y_0) \leq \epsilon]$ equals one if $d(y, y_0) \leq \epsilon$ and zero otherwise. Importantly, while $p(y|\theta)$ is generally unknown for simulator-based models, it is still possible to use $p_{d,\epsilon}(\theta, y|y_0)$ as the target distribution in a Metropolis-Hastings MCMC algorithm by choosing the proposal distribution in the right way. The obtained (marginal) samples of θ then follow the approximate posterior $p_{d,\epsilon}(\theta|y_0)$.

Assuming that the Markov chain is at iteration i in state $x^{(i)} = (\theta^{(i)}, y^{(i)})$ where $d(y^{(i)}, y_0) \leq \epsilon$, the Metropolis-Hastings algorithm involves sampling candidate states $x = (\theta, y)$ from a proposal distribution $q(x|x^{(i)})$ and accepting the candidates with the probability $A(x|x^{(i)})$,

$$A(x|x^{(i)}) = \min \left(1, \frac{p_{d,\epsilon}(x|y_0)q(x^{(i)}|x)}{p_{d,\epsilon}(x^{(i)}|y_0)q(x|x^{(i)})} \right). \quad (9)$$

Choosing the proposal distribution such that the move from $x^{(i)} = (\theta^{(i)}, y^{(i)})$ to $x = (\theta, y)$ does not depend on the value of $y^{(i)}$, and that y is sampled from the simulator-based model with parameter value θ irrespective of $\theta^{(i)}$, we have

$$q(x|x^{(i)}) = q(\theta|\theta^{(i)})p(y|\theta), \quad (10)$$

where $q(\theta|\theta^{(i)})$ is a suitable proposal distribution for θ . As a result of this choice, the unknown quantities in Equation (9) cancel out,

$$\begin{aligned} A(x|x^{(i)}) &= \min \left(1, \frac{p(\theta)}{p(\theta^{(i)})} \frac{p(y|\theta)}{p(y^{(i)}|\theta^{(i)})} \frac{\mathbb{1}[d(y, y_0) \leq \epsilon]}{\mathbb{1}[d(y^{(i)}, y_0) \leq \epsilon]} \frac{q(\theta^{(i)}|\theta)}{q(\theta|\theta^{(i)})} \frac{p(y^{(i)}|\theta^{(i)})}{p(y|\theta)} \right) \\ &= \min \left(1, \frac{p(\theta)}{p(\theta^{(i)})} \frac{q(\theta^{(i)}|\theta)}{q(\theta|\theta^{(i)})} \frac{\mathbb{1}[d(y, y_0) \leq \epsilon]}{\mathbb{1}[d(y^{(i)}, y_0) \leq \epsilon]} \right) \\ &= \mathbb{1}[d(y, y_0) \leq \epsilon] \min \left(1, \frac{p(\theta)}{p(\theta^{(i)})} \frac{q(\theta^{(i)}|\theta)}{q(\theta|\theta^{(i)})} \right). \end{aligned} \quad (11)$$

This means that the acceptance probability is only probabilistic in θ since a proposal (θ, y) is immediately rejected if the condition $d(y, y_0) \leq \epsilon$ is not met. While the Markov chain operates in the (θ, y) space, the choice of the proposal distribution decouples the acceptance criterion into an ordinary Metropolis-Hastings criterion for θ and the previously seen ABC rejection criterion for y . The resulting algorithm, shown in full in the appendix, is known as MCMC ABC algorithm and was introduced by Marjoram et al. (2003).

An advantage of the MCMC ABC algorithm is that the parameter values do not need to be drawn from the prior, which most often hampers the rejection sampler by incurring a high rejection rate of the proposals. As the Markov chain converges, the proposed parameter values follow the posterior with some added noise. A potential disadvantage, however, is the continuing presence of the rejection condition $d(y, y_0) \leq \epsilon$ which dominates the acceptance rate of the algorithm. Parameters in the tails of the posteriors have, by definition, a small probability to generate data y_θ satisfying the rejection condition, which can lead to a “sticky” Markov chain where the state tends to remain constant for many iterations.

Sequential Monte Carlo ABC

The sequential Monte Carlo ABC algorithm can be considered as an adaptation of importance sampling which is a popular technique in statistics (see, for example, Robert and Casella 2004, Chapter 3). If one uses a general distribution $\phi(\theta)$ in place of the prior $p(\theta)$, Algorithm 2 produces samples that follow a distribution proportional to $\phi(\theta) \Pr(d(Y_\theta, y_0) \leq \epsilon)$. However, by weighting the accepted parameters $\theta^{(i)}$ with $w^{(i)}$,

$$w^{(i)} \propto \frac{p(\theta^{(i)})}{\phi(\theta^{(i)})}, \quad (12)$$

the resulting weighted samples follow $p_{d,\epsilon}(\theta|y_0)$. This kind of trick is used in importance

sampling and can be employed in ABC to iteratively morph the prior into a posterior.

The basic idea is to use a sequence of shrinking thresholds ϵ_t and to define the proposal distribution ϕ_t at iteration t based on the weighted samples $\theta_{t-1}^{(i)}$ from the previous iteration (Figure 8). This is typically done by defining a mixture distribution,

$$\phi_t(\theta) = \frac{1}{N} \sum_{i=1}^N q_t(\theta|\theta_{t-1}^{(i)}) w_{t-1}^{(i)}, \quad (13)$$

where $q_t(\theta|\theta_{t-1}^{(i)})$ is often a Gaussian distribution with mean $\theta_{t-1}^{(i)}$ and a covariance matrix estimated from the samples. Sampling from ϕ_t can be done by choosing $\theta_{t-1}^{(i)}$ with probability $w_{t-1}^{(i)}$ and then perturbing the chosen parameter according to q_t . The proposed sample is then accepted or rejected as in Algorithm 2 and the weights of the accepted samples are computed with Equation 12. Such iterative algorithms were proposed by Sisson et al. (2007); Beaumont et al. (2009); Toni et al. (2009) and are called Sequential Monte Carlo (SMC) ABC algorithms or Population Monte Carlo (PMC) ABC algorithms. The algorithm by Beaumont et al. (2009) is given in the appendix.

Similar to the MCMC ABC, the samples proposed by the SMC algorithm follow the posterior $p_{d,\epsilon_t}(\theta|y_0)$ with some added noise. The proposed parameter values are drawn from the prior only at the first iteration after which adaptive proposal distributions ϕ_t closer to the true posterior are used (see Figure 8 for an illustration). This reduces the running time as the number of rejections is lower compared to the basic rejection ABC algorithm. For small values of ϵ , however, the probability to accept a parameter value becomes very small, even if the parameter value was sampled from the true posterior. This results in long computing times in the final iterations of the algorithm without notable improvements in the approximation of the posterior.

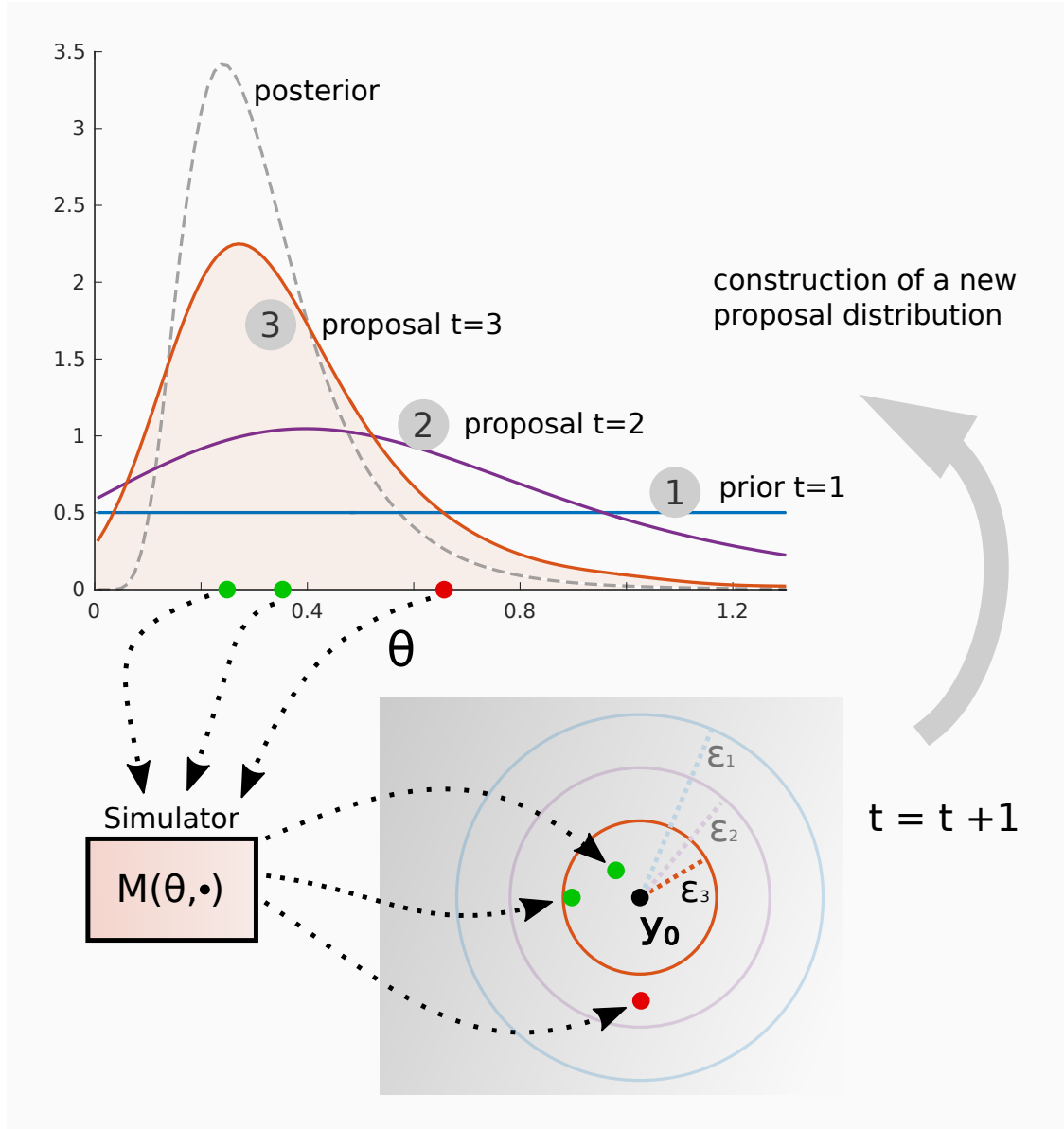


Figure 8: Illustration of sequential Monte Carlo ABC using the tuberculosis example. The first proposal distribution is the prior and the threshold value used is ϵ_1 . The proposal distribution in iteration t is based on the sample of size N from the previous iteration. The threshold value ϵ_t is decreased at every iteration as the proposal distributions become similar to the true posterior. The figure shows parameters drawn from the proposal distribution of the third iteration ($t = 3$). The red proposal is rejected because the corresponding simulation outcome is too far from the observed data y_0 . At iteration $t = 2$, however, it would have been accepted. After iteration t , the accepted parameter values follow the approximate posterior $p_{d, \epsilon_t}(\theta|y_0)$. As long as the threshold values ϵ_t decrease, the approximation becomes more accurate at each iteration.

Post-sampling correction methods

We assume here that the distance $d(y_\theta, y_0)$ is specified in terms of summary statistics, that is, $d(y_\theta, y_0) = d_T(t_\theta, t_0)$, with $t_\theta = T(y_\theta)$ and $t_0 = T(y_0)$. As ϵ decreases to zero, the approximate posterior $p_{d,\epsilon}(\theta|y_0)$ in Equation (5) converges to $p(\theta|t_0)$, where we use $p(\theta|t)$ to denote the conditional distribution of θ given a value of the summary statistics t . While small values of ϵ are thus preferred in theory, making them too small is not feasible in practice because of the correspondingly small acceptance rate and the resulting large Monte Carlo error. We here present two schemes that aim at adjusting $p_{d,\epsilon}(\theta|y_0)$ without further sampling so that the adjusted distribution is closer to $p(\theta|t_0)$.

For the first scheme, we note that if we had a mechanism to sample from $p(\theta|t)$, we could sample from the limiting approximate posterior by using $t = t_0$. The post-sampling correction methods in the first scheme thus estimate $p(\theta|t)$ and use the estimated conditional distributions to sample from $p(\theta|t_0)$. In order to facilitate sampling, $p(\theta|t)$ is expressed in terms of a generative (regression) model,

$$\theta = f(t, \xi), \tag{14}$$

where f is a vector-valued function and ξ a vector of random variables for residuals. By suitably defining f , we can assume that the random variables of the vector ξ are independent, of zero mean and equal variance, and that their distribution p_ξ does not depend on t . Importantly, the model does not need to hold for all t because, ultimately, we would like to sample from it using $t = t_0$ only. Assuming that the model f holds for $d_T(t, t_0) \leq \delta$ and that we have (weighted) samples $(t^{(i)}, \tilde{\theta}^{(i)}) = (T(y_\theta^{(i)}), \tilde{\theta}^{(i)})$ available from an ABC algorithm with a threshold $\epsilon \leq \delta$, the model f can be estimated by regressing θ on the summary statistics t .

In order to sample θ using the estimated model \hat{f} , we need to know the distribution

of ξ . For that, the residuals $\xi^{(i)}$ are determined by solving the regression equation,

$$\tilde{\theta}^{(i)} = \hat{f}(t^{(i)}, \xi^{(i)}). \quad (15)$$

The residuals $\xi^{(i)}$ can be used to estimate p_ξ , or as usually is the case in ABC, be directly employed in the sampling of the θ ,

$$\theta^{(i)} = \hat{f}(t_0, \xi^{(i)}). \quad (16)$$

If the original samples $(t^{(i)}, \tilde{\theta}^{(i)})$ are weighted, both the $\xi^{(i)}$ and the new “adjusted” samples $\theta^{(i)}$ inherit the weights. By construction, if the relation between t and θ is estimated correctly, the (weighted) samples $\theta^{(i)}$ follow $p_{d,\epsilon}(\theta|y_0)$ with $\epsilon = 0$.

In most models f employed so far, the individual components of θ are treated separately, thus not accounting for possible correlations between them. For this paragraph we thus let θ be a scalar. The first regression model used was linear (Beaumont et al. 2002),

$$\theta = f_1(t, \xi), \quad f_1(t, \xi) = \alpha + (t - t_0)^\top \beta + \xi, \quad (17)$$

which results in the adjustment $\theta^{(i)} = \tilde{\theta}^{(i)} - (t^{(i)} - t_0)^\top \hat{\beta}$ where $\hat{\beta}$ is the learned regression coefficient (Figure 9). When applied to the model of the spread of tuberculosis, with summary statistic T_1 (see Equation 7), the adjustment is able to correct the bias caused by the non-zero threshold $\epsilon = 0.1$, e.g. the estimated model \hat{f} is accurate (Figure 10). With summary statistic T_2 the threshold $\epsilon = 0.1$ is too large for accurate adjustment, although the result is still closer to the target distribution than the original. Note also that here the target distribution of the adjustment is substantially different from the true posterior due to the bias incurred by summary statistic T_2 .

Also non-linear models f have been proposed. Blum (2010) assumed a quadratic

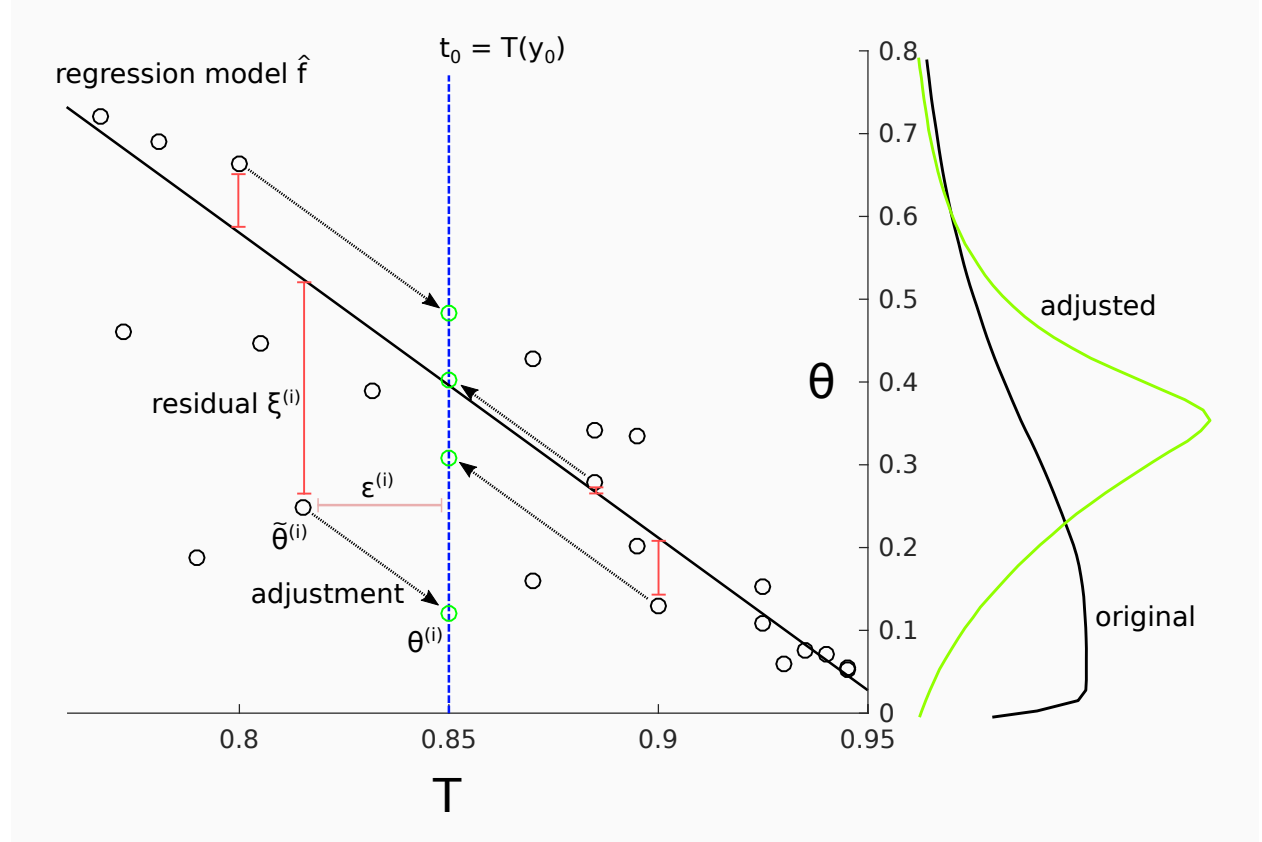


Figure 9: Illustration of the linear regression adjustment (Beaumont et al. 2002). First the regression model \hat{f} is learned and then, based on \hat{f} , the simulations are adjusted as if they were sampled from $p_{d,\epsilon}(\theta|y_0)$ with $\epsilon = 0$. Note that the residuals $\xi^{(i)}$ are preserved. The change in the posterior densities after the adjustment is shown on the right. Here the black (original) and green (adjusted) curves correspond to the respectively named curves in Figure 10(b).

model,

$$\theta = f_2(t, \xi), \quad f_2(t, \xi) = \alpha + (t - t_0)^\top \beta + \frac{1}{2}(t - t_0)^\top \gamma (t - t_0) + \xi, \quad (18)$$

where γ is a symmetric matrix, that adds a quadratic term to the linear adjustment. A more general nonlinear model was considered by Blum and François (2010),

$$\theta = f_3(t, \xi), \quad f_3(t, \xi) = m(t) + \sigma(t)\xi, \quad (19)$$

where $m(t)$ models the conditional mean and $\sigma(t)$ the conditional standard deviation of θ . Both functions were fitted using a multi-layer neural network, and denoting the learned functions by \hat{m} and $\hat{\sigma}$, the following adjustments were obtained

$$\theta^{(i)} = \hat{m}(t_0) + \hat{\sigma}(t_0)\hat{\sigma}(t^{(i)})^{-1}(\tilde{\theta}^{(i)} - \hat{m}(t^{(i)})). \quad (20)$$

The term $\hat{m}(t_0)$ is an estimate of the posterior mean of θ while $\hat{\sigma}(t_0)$ is an estimate of the posterior standard deviation of the parameter. They can both be used to succinctly summarize the posterior distribution of θ .

A more complicated model $f(t, \xi)$ is not necessarily better than a simpler one. It depends on the amount of the training data available to fit it, that is, the amount of original samples $(t^{(i)}, \tilde{\theta}^{(i)})$ that satisfy $d_T(t, t_0) \leq \delta$. The different models presented above were compared by Blum and François (2010) who also pointed out that techniques for model selection from the regression literature can be used to select among them.

While the first scheme to adjust $p_{d,\epsilon}(\theta|y_0)$ consists of estimating $p(\theta|t)$, the second scheme consists of estimating $p(t|\theta)$, that is the conditional distribution of the summary statistics given a parameter value. The rationale of this approach is that knowing $p(t|\theta)$ implies knowing the approximate likelihood function $L_{d,\epsilon}(\theta)$ for $\epsilon \rightarrow 0$, because

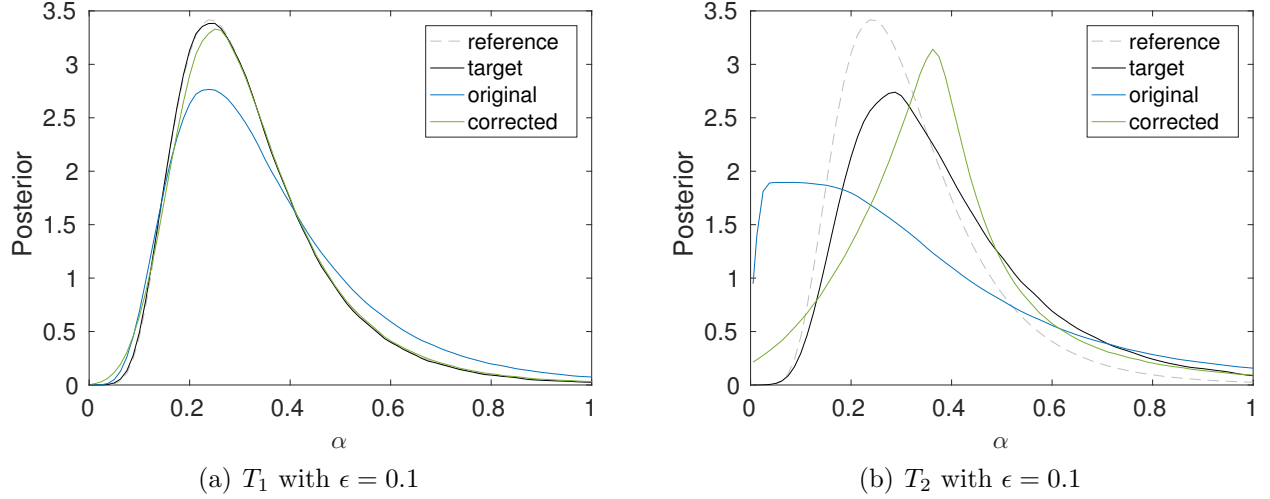


Figure 10: Linear regression adjustment (Beaumont et al. 2002) applied to the example model of the spread of tuberculosis (compare to Figure 5). The target distribution of the adjustment is the posterior $p_{d,\epsilon}(\theta|y_0)$ with the threshold decreased to $\epsilon = 0$. Note that when using summary statistic T_2 the target distribution is substantially different from the true posterior (reference) because of the bias incurred by T_2 .

$p(t_0|\theta) = \lim_{\epsilon \rightarrow 0} L_{d,\epsilon}(\theta)$ when the distance $d(y_\theta, y_0)$ is specified in terms of summary statistics.

Importantly, $p(t|\theta)$ does not need to be known everywhere but only locally around t_0 , where $d_T(t, t_0) \leq \epsilon$. If we use $p_\epsilon(t|\theta)$ to denote the distribution of t conditional on θ and $d_T(t, t_0) \leq \epsilon$, Leuenberger and Wegmann (2010) showed that $p_\epsilon(t_0|\theta)$ takes the role of a local likelihood function and $p_{d,\epsilon}(\theta|y_0)$ the role of a local prior, and that the local posterior equals the true posterior $p(\theta|t_0)$.

The functional form of $p_\epsilon(t|\theta)$ is generally not known. However, as in the first scheme, running an ABC algorithm with threshold ϵ provides data $(t^{(i)}, \tilde{\theta}^{(i)})$ that can be used to estimate a model of $p_\epsilon(t|\theta)$. Since the model does not need to hold for all values of the summary statistics, but only for those in the neighborhood of t_0 , Leuenberger and Wegmann (2010) proposed to model $p_\epsilon(t|\theta)$ as Gaussian with constant covariance matrix and a mean depending linearly on θ . When the samples $(t^{(i)}, \tilde{\theta}^{(i)})$ are used to approximate

$p_{d,\epsilon}(\theta|y_0)$ as a kernel-density estimate, the Gaussianity assumption on $p_\epsilon(t|\theta)$ facilitates the derivation of closed-form formulae to adjust the kernel-density representation of $p_{d,\epsilon}(\theta|y_0)$ so that it becomes an approximation of $p(\theta|t_0)$ (Leuenberger and Wegmann 2010).

While Leuenberger and Wegmann (2010) modeled $p_\epsilon(t|\theta)$ as Gaussian, other models can be used as well. Alternatively, one may make the mean of the Gaussian depend nonlinearly on θ and allow the covariance of the summary statistic depend on θ . This was done by Wood (2010) and the model was found rich enough to represent $p(t|\theta)$ for all values of the summary statistics and not only for those in the neighborhood of the observed one.

RECENT DEVELOPMENTS

We here present recent advances that aim to make approximate Bayesian computation both computationally and statistically more efficient. This presentation focuses on our own work (Gutmann et al. 2014; Gutmann and Corander 2015).

Computational efficiency

The computational cost of ABC can be attributed to two main factors:

1. Most of the parameter values result in large distances between the simulated and observed data and those parameter values for which the distances tend to be small are unknown.
2. Generating simulated data sets, that is, running the simulator, may be costly.

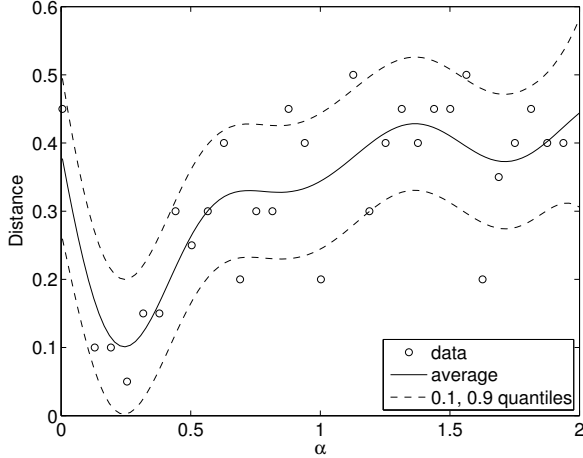
MCMC ABC and SMC ABC were partly introduced to avoid proposing parameters in regions where the distance is large. Nonetheless, typically millions of simulations are needed to infer the posterior distribution of a handful of parameters only. A key obstacle to efficiency in these algorithms is the continued presence of the rejection mechanism

$d(y_\theta, y_0) \leq \epsilon$, or more generally, the online decisions about the similarity between y_θ and y_0 . In recent work, Gutmann and Corander (2015) proposed a framework called Bayesian optimization for likelihood-free inference (BOLFI) for performing approximate Bayesian computation which overcomes this obstacle by learning a probabilistic model about the stochastic relation between the parameter values and the distance (Figure 11). After learning, the model can be used to approximate $L_{d,\epsilon}(\theta)$, and thus $p_{d,\epsilon}(\theta|y_0)$, for any ϵ without requiring further runs of the simulator (Figure 12).

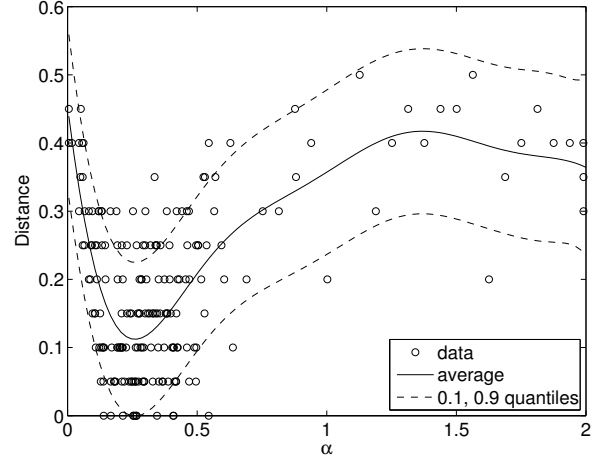
Like the post-sampling correction methods presented in the previous section, BOLFI relies on a probabilistic model to make ABC more efficient. However, the quantities modeled differ, since in the post-sampling correction methods the relation between summary statistics and parameters is modeled, while BOLFI focuses on the relation between the parameters and the distance. A potential advantage of the latter approach is that the distance is a univariate quantity while the parameters and summary statistics may be multi-dimensional. Furthermore, BOLFI does not assume that the distance is defined via summary statistics, and can be used without first running another ABC algorithm.

Learning of the model of $d(Y_\theta, y_0)$ requires data about the relation between θ and $d(Y_\theta, y_0)$. In BOLFI, the data are actively acquired focusing on regions of the parameter space where the distance tends to be small. This is achieved by leveraging techniques from Bayesian optimization (see for example Jones 2001; Brochu et al. 2010), hence its name. Ultimately, the framework provided by Gutmann and Corander (2015) reduces the computational cost of ABC by addressing both of the factors mentioned above. The first point is addressed by learning from data which parameter values tend to have small distances, while the second problem is resolved by focusing on areas where the distance tends to be small when learning the model and by not requiring further runs of the simulator once the model is learned.

While BOLFI is not restricted to a particular model for $d(Y_\theta, y_0)$, Gutmann and



(a) After initialization (30 data points)



(b) After active data acquisition (200 data points)

Figure 11: The basic idea of Bayesian optimization in likelihood-free inference (BOLFI) is to model the distance, and to prioritize regions of the parameter space where the distance tends to be small. The solid curves show the modeled average behavior of the distance $d_1(Y_\theta, y_0)$, and the dashed curves its variability for the tuberculosis example.

Corander (2015) used Gaussian processes in the applications in their paper. Gaussian processes have also been used in other work as surrogate models for quantities that are expensive to compute. Wilkinson (2014) used them to model the logarithm of $L_{d,\epsilon}(\theta)$, and the training data were constructed based on quasi-random numbers covering the parameter space. Meeds and Welling (2014) used Gaussian processes to model the empirical mean and covariances of the summary statistics as a function of θ . Instead of simulating these quantities for every θ , values from the model were used in a Markov chain Monte Carlo algorithm in approximating the likelihood. These approaches have been demonstrated to assist in speeding up approximate Bayesian computation.

Statistical efficiency

We have seen that the statistical efficiency of ABC algorithms depends heavily on the summary statistics chosen, the distance between them, and the locality of the inference. In

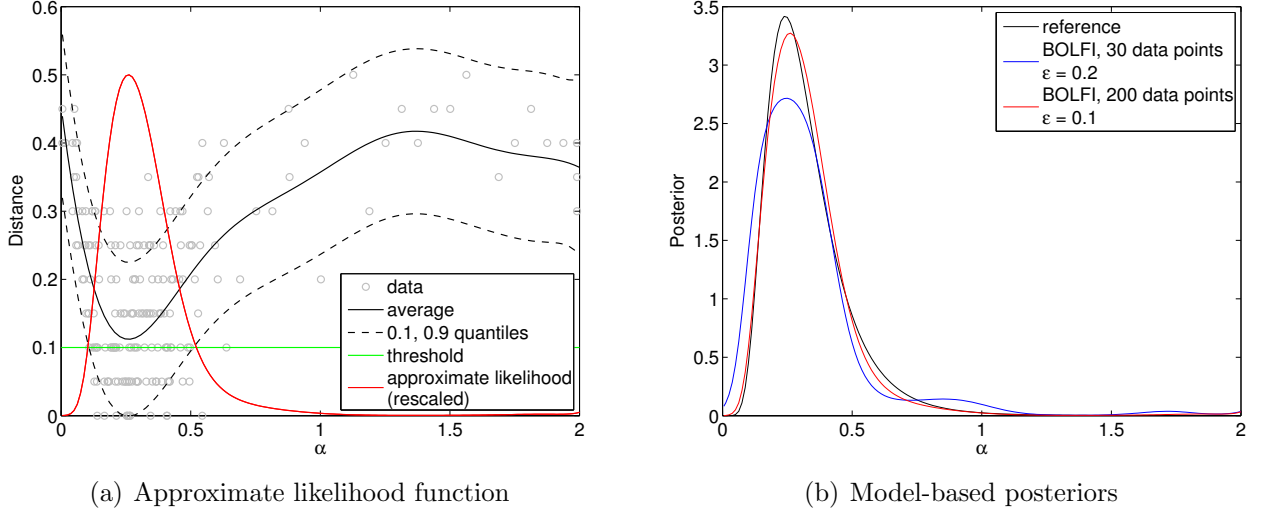


Figure 12: In BOLFI, the estimated model of $d(y_\theta, y_0)$ is used to approximate $L_{d,\epsilon}(\theta)$ by computing the probability that the distance is below a threshold ϵ . This kind of likelihood approximation leads to a model-based approximation of $p_{d,\epsilon}(\theta|y_0)$. The KL-divergence between the reference solution and the BOLFI solution with 30 data points is 0.09, and for 200 data points it is 0.01. Comparison with Figure 6 shows that BOLFI increases the computational efficiency of ABC by several orders of magnitude.

a recent work, Gutmann et al. (2014) formulated the problem of measuring the distance between simulated and observed data as a classification problem: Two data sets are judged maximally similar if they cannot be told apart significantly above chance level (50% accuracy in the classification problem). On the other hand, two data sets are maximally dissimilar if they can be told apart with 100% classification accuracy. In essence, classification is used to assess the distance between simulated and observed data.

The classification rule used to measure the distance was learned from the data, which simplifies the inference since only a function (hypothesis) space needs to be pre-specified by the user. In the process, Gutmann et al. (2014) also chose a subset or weighted (nonlinear) combination of summary statistics to achieve the best classification accuracy. This choice depended on the parameter values used to generate the simulated data. While computationally more expensive than the traditional approach, the classifier

approach has the advantage of being a data-driven way to measure the distance between the simulated and observed data which respects the locality of the inference.

VALIDATION OF APPROXIMATE BAYESIAN COMPUTATION

Due to the several levels of approximation, it is generally a recommendable practice to perform validity analyses of the ABC inferences. We here discuss some of the possibilities suggested in the literature.

The ability to generate data from simulator-based models enables basic sanity checks for the feasibility of the inference with a given setting and algorithm. The general approach is to perform inference where synthetic data sets y_0^* are generated with known parameter values θ^* to play the role of the observed data y_0 . To assess whether the posterior distribution is concentrated around the right parameter values, one may then compute the average error between the posterior mean (mode) and θ^* , or the expected squared distance between the posterior samples and θ^* (Wegmann et al. 2009). To assess whether the spread of the posterior distribution is not overly large or small, one may compute confidence (credibility) intervals and check their coverage. When the nominal confidence levels are accurate, 95% confidence intervals, for example, should contain θ^* in 95% of the simulation experiments (Wegmann et al. 2009; Prangle et al. 2014). Such tests can be performed *a priori*, by sampling y_0^* from the prior before having seen the actual data to be analyzed, or also *a posteriori*, by sampling y_0^* from the inferred posterior or from the prior restricted to some area of interest (Prangle et al. 2014). Corresponding techniques have also been suggested for the purpose of specifying the threshold value ϵ as discussed earlier in this paper. It can be also beneficial here to store the generated data sets together with their parameter values so that the validations can be run without having to re-generate new data on every occasion.

The ABC framework provides a straightforward way to investigate the goodness-of-fit of the model. The distances $d(y_\theta, y_0)$ indicate how close the simulated data y_θ are to the observed data y_0 . If all of the distances remain large, it may be an indication of a deficient model, as the model is unable to produce data similar to the observed data. Ratmann et al. (2009) proposed a method called ABC under model uncertainty (ABC μ) where they augment the likelihood with unknown error terms for each of the different summary statistics used. The error terms are assumed to have mean zero and are sampled together with the parameters of the model. If however the mean of the error terms is found to deviate from 0, it may indicate a systematic error in the model.

Yet another issue is to consider identifiability of the model given the observed data. The likelihood function indicates the extent to which parameter values are congruent with the observed data. A strong curvature at its maximum indicates that the maximizing parameter value is clearly to be preferred while a minor curvature means that several other parameter values are nearly equally supported by the data. More generally, if the likelihood surface is mostly flat over the parameter space, the data are not providing sufficient information to identify the model parameters. While the likelihood function is generally not available for simulator-based models, the arguments provided do also hold for the approximate likelihood function $L_{d,\epsilon}(\theta)$ in Equation (6). On one hand, the approximate likelihood function can be used to investigate the identifiability of the simulator-based model. On the other hand, it allows one to assess the quality of the distance d or threshold ϵ chosen. Flat approximate likelihood surfaces, for instance, indicate that ϵ could be too large or that the distance function d is not able to accurately measure differences between the data sets.

The approximate likelihood $L_{d,\epsilon}(\theta)$ can be obtained either by the method of Gutmann and Corander (2015) or also by any other ABC algorithm by assuming a uniform prior on a region of interest. Lintusaari et al. (2016) used such an approach to investigate

the identifiability of the tuberculosis model considered as example in the previous sections, and to compare different distance functions. Further, one may (visually) compare the (marginal) prior and the inferred (marginal) posterior (e.g. Blum 2010). Both approaches are applicable not only to the real observed data y_0 but also to the synthetic data y_0^* for which the data-generating parameters θ^* are known. If the employed ABC algorithm is working appropriately, both $L_{d,\epsilon}(\theta)$ and the posteriors should clearly change when the characteristics of the observed data change markedly. In particular, if the number of observations is increased, the approximate likelihood and posterior should in general become more concentrated around the data-generating parameter values. While failure to pass such sanity checks may be an indicator that the choice of d and ϵ could be improved, it can also indicate that the model may not be fully identifiable.

CONCLUSIONS

It is possible to model complex biological phenomena in a realistic manner with the aid of simulator-based models. However the likelihood function for such models is usually intractable and raises serious methodological challenges to perform statistical inference. Approximate Bayesian computation has become synonymous for approximate Bayesian inference for simulator-based models. We have here reviewed its foundations, the most widely considered inference algorithms, together with recent advances that increase its statistical and computational efficiency.

While the review is solely restricted to Bayesian methods, there exists a large body of literature on non-Bayesian approaches, for instance, the methods of simulated moments (Pakes and Pollard 1989; McFadden 1989) or indirect inference (Gouriéroux et al. 1993; Heggland and Frigessi 2004), both having their origin in econometrics.

We focused on the central topics related to parameter inference with ABC. Nevertheless, ABC is also applicable to model selection (see, for example, the review by

Marin et al. 2012), and while we have reviewed methods making the basic ABC algorithms more efficient, we have not discussed the important topic of how to use ABC for high-dimensional inference. We point the interested readers to the work by Li et al. (2015) and also to the discussion by Gutmann and Corander (2015).

For practical purpose, there exist multiple software packages implementing the different ABC algorithms, summary statistic selection, validation methods, post processing, and ABC model selection methods. Nunes and Prangle (2015) provide a recent list of available packages with information about their implementation language, platform and targeted field of study. In summary, approximate Bayesian computation is currently a very active methodological research field, and this activity will likely result in several advances to improve its applicability to answering important biological research questions in the near future.

ACKNOWLEDGEMENTS

We acknowledge the computational resources provided by the Aalto Science-IT project.

FUNDING

This work was supported by the Academy of Finland (Finnish Centre of Excellence in Computational Inference Research COIN).

Author contributions: JL, MUG, RD and JC wrote the paper; JL performed the simulations; SK contributed to writing and planning of the paper.

*

References

- Aeschbacher, S., M. Beaumont, and A. Futschik. 2012. A novel approach for choosing summary statistics in approximate Bayesian computation. *Genetics* 192:1027–1047.
- Anderson, R. M. and R. M. May. 1992. *Infectious Diseases of Humans: Dynamics and Control*. Oxford University Press.
- Barber, S., J. Voss, and M. Webster. 2015. The rate of convergence for approximate Bayesian computation. *Electronic Journal of Statistics* Pages 80–105.
- Baudet, C., B. Donati, B. Sinimeri, P. Crescenzi, C. Gautier, C. Matias, and M.-F. Sagot. 2015. Cophylogeny reconstruction via an approximate bayesian computation. *Syst Biol* 64:416–431.
- Beaumont, M. A. 2010. Approximate Bayesian computation in evolution and ecology. *Annual review of ecology, evolution, and systematics* 41:379–406.
- Beaumont, M. A., J.-M. Cornuet, J.-M. Marin, and C. P. Robert. 2009. Adaptive approximate Bayesian computation. *Biometrika* 96:983–990.
- Beaumont, M. A., W. Zhang, and D. J. Balding. 2002. Approximate Bayesian computation in population genetics. *Genetics* 162:2025–2035.
- Biau, G., F. Cérou, and A. Guyader. 2015. New insights into approximate Bayesian computation. *Annales de l’Institut Henri Poincaré (B) Probabilités et Statistiques* 51:376–403.
- Blum, M. and O. François. 2010. Non-linear regression models for approximate Bayesian computation. *Statistics and Computing* 20:63–73.
- Blum, M. G. B. 2010. Approximate Bayesian computation: A nonparametric perspective. *Journal of the American Statistical Association* 105:1178–1187.

- Blum, M. G. B., M. A. Nunes, D. Prangle, and S. A. Sisson. 2013. A comparative review of dimension reduction methods in approximate Bayesian computation. *Statistical Science* 28:189–208.
- Brochu, E., V. Cora, and N. de Freitas. 2010. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *arXiv:1012.2599* .
- Curat, M. and L. Excoffier. 2004. Modern humans did not admix with neanderthals during their range expansion into europe. *PLoS Biol* 2:e421.
- Diggle, P. J. and R. J. Gratton. 1984. Monte carlo methods of inference for implicit statistical models. *Journal of the Royal Statistical Society. Series B (Methodological)* 46:193–227.
- Excoffier, L., I. Dupanloup, E. Huerta-Sánchez, V. C. Sousa, and M. Foll. 2013. Robust demographic inference from genomic and snp data. *PLoS Genet* 9:e1003905.
- Fagundes, N. J. R., N. Ray, M. Beaumont, S. Neuenschwander, F. M. Salzano, S. L. Bonatto, and L. Excoffier. 2007. Statistical evaluation of alternative models of human evolution. *Proceedings of the National Academy of Sciences* 104:17614–17619.
- Faisal, M., A. Futschik, and I. Hussain. 2013. A new approach to choose acceptance cutoff for approximate Bayesian computation. *Journal of Applied Statistics* 40:862–869.
- Fan, H. H. and L. S. Kubatko. 2011. Estimating species trees using approximate bayesian computation. *Molecular Phylogenetics and Evolution* 59:354 – 363.
- Fearnhead, P. and D. Prangle. 2012. Constructing summary statistics for approximate Bayesian computation: semi-automatic approximate Bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 74:419–474.

- Gouriéroux, C., A. Monfort, and E. Renault. 1993. Indirect inference. *Journal of Applied Econometrics* 8:S85–S118.
- Green, P., K. Latuszynski, M. Pereyra, and C. P. Robert. 2015. Bayesian computation: a summary of the current state, and samples backwards and forwards. *Statistics and Computing* 25:835–862.
- Gutmann, M. and J. Corander. 2015. Bayesian optimization for likelihood-free inference of simulator-based statistical models. *Journal of Machine Learning Research* in press.
- Gutmann, M., R. Dutta, S. Kaski, and J. Corander. 2014. Statistical inference of intractable generative models via classification. *arXiv:1407.4981* .
- Heggland, K. and A. Frigessi. 2004. Estimating functions in indirect inference. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 66:447–462.
- Hickerson, M. J., E. A. Stahl, and H. A. Lessios. 2006. Test for simultaneous divergence using approximate bayesian computation. *Evolution* 60:2435–2453.
- Itan, Y., A. Powell, M. A. Beaumont, J. Burger, and M. G. Thomas. 2009. The origins of lactase persistence in europe. *PLoS Comput Biol* 5:e1000491.
- Jones, D. 2001. A taxonomy of global optimization methods based on response surfaces. *Journal of Global Optimization* 21:345–383.
- Joyce, P. and P. Marjoram. 2008. Approximately sufficient statistics and Bayesian computation. *Statistical Applications in Genetics and Molecular Biology* 26.
- Kullback, S. and R. A. Leibler. 1951. On information and sufficiency. *Ann. Math. Statist.* 22:79–86.

- Leuenberger, C. and D. Wegmann. 2010. Bayesian computation and model selection without likelihoods. *Genetics* 184:243–252.
- Li, J., D. J. Nott, and S. A. Sisson. 2015. Extending approximate Bayesian computation methods to high dimensions via Gaussian copula. *arXiv:1504.04093* .
- Lintusaari, J., M. U. Gutmann, S. Kaski, and J. Corander. 2016. On the identifiability of transmission dynamic models for infectious diseases. *Genetics* .
- Marin, J.-M., P. Pudlo, C. Robert, and R. Ryder. 2012. Approximate Bayesian computational methods. *Statistics and Computing* 22:1167–1180.
- Marjoram, P., J. Molitor, V. Plagnol, and S. Tavaré. 2003. Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences* 100:15324–15328.
- Marttinen, P., N. J. Croucher, M. U. Gutmann, J. Corander, and W. P. Hanage. 2015. Recombination produces coherent bacterial species clusters in both core and accessory genomes. *Microbial Genomics*. Published ahead of print.
- McFadden, D. 1989. A method of simulated moments for estimation of discrete response models without numerical integration. *Econometrica* 57:995–1026.
- Meeds, E. and M. Welling. 2014. GPS-ABC: Gaussian process surrogate approximate Bayesian computation. *in* Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence (UAI).
- Nunes, M. A. and D. J. Balding. 2010. On optimal selection of summary statistics for approximate Bayesian computation. *Statistical Applications in Genetics and Molecular Biology* 9.
- Nunes, M. A. and D. Prangle. 2015. abctools: An R package for tuning approximate Bayesian computation analyses. *The R Journal*. To appear.

- Pakes, A. and D. Pollard. 1989. Simulation and the asymptotics of optimization estimators. *Econometrica* 57:1027–1057.
- Prangle, D., M. G. B. Blum, G. Popovic, and S. A. Sisson. 2014. Diagnostic tools for approximate Bayesian computation using the coverage property. *Australian & New Zealand Journal of Statistics* 56:309–329.
- Pritchard, J. K., M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman. 1999. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Molecular Biology and Evolution* 16:1791–1798.
- Ratmann, O., C. Andrieu, C. Wiuf, and S. Richardson. 2009. Model criticism based on likelihood-free inference, with an application to protein network evolution. *Proceedings of the National Academy of Sciences* 106:10576–10581.
- Ratmann, O., G. Donker, A. Meijer, C. Fraser, and K. Koelle. 2012. Phylodynamic inference and model assessment with approximate bayesian computation: influenza as a case study. *PLoS Comput Biol* 8:e1002835.
- Robert, C. and G. Casella. 2004. *Monte Carlo Statistical Methods*. 2 ed. Springer.
- Rubin, D. B. 1984. Bayesianly justifiable and relevant frequency calculations for the applied statistician. *The Annals of Statistics* 12:1151–1172.
- Silk, D., S. Filippi, and M. P. Stumpf. 2013. Optimizing threshold-schedules for sequential approximate bayesian computation: applications to molecular systems. *Statistical applications in genetics and molecular biology* 12:603–618.
- Sisson, S. A., Y. Fan, and M. M. Tanaka. 2007. Sequential Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences* 104:1760–1765.

- Slater, G. J., L. J. Harmon, D. Wegmann, P. Joyce, L. J. Revell, and M. E. Alfaro. 2012. Fitting models of continuous trait evolution to incompletely sampled comparative data using approximate Bayesian computation. *Evolution* 66:752–762.
- Stadler, T. 2011. Inferring epidemiological parameters on the basis of allele frequencies. *Genetics* 188:663–672.
- Tanaka, M. M., A. R. Francis, F. Luciani, and S. A. Sisson. 2006. Using approximate Bayesian computation to estimate tuberculosis transmission parameters from genotype data. *Genetics* 173:1511–1520.
- Toni, T., D. Welch, N. Strelkowa, A. Ipsen, and M. P. Stumpf. 2009. Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of The Royal Society Interface* 6:187–202.
- Wegmann, D., C. Leuenberger, and L. Excoffier. 2009. Efficient approximate Bayesian computation coupled with Markov chain Monte Carlo without likelihood. *Genetics* 182:129–141.
- Wilkinson, R. 2014. Accelerating ABC methods using Gaussian processes. *in* Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS).
- Wood, S. 2010. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature* 466:1102–1104.

APPENDIX

For completeness, we state below the algorithms for MCMC-ABC and SMC-ABC by Marjoram et al. (2003) and Beaumont et al. (2009), respectively.

Algorithm 3 MCMC-ABC algorithm producing N samples from the approximate posterior distribution $p_{d,\epsilon}(\theta|y_0)$

Require: Set the initial value $\theta^{(0)}$

```
1: for  $i = 1$  to  $N$  do
2:   Generate  $\theta$  from a transition kernel  $q(\cdot|\theta^{(i-1)})$ 
3:   Generate  $y_\theta$  from the simulator
4:   if  $d(y_\theta, y_0) \leq \epsilon$  then
5:     Calculate  $A = A(\theta|\theta^{(i-1)}) = p(\theta)q(\theta^{(i-1)}|\theta)/(p(\theta^{(i-1)})q(\theta|\theta^{(i-1)}))$ 
6:     Generate  $u$  from  $\text{Uni}(0, 1)$ 
7:     if  $u < A$  then
8:        $\theta^{(i)} \leftarrow \theta$ 
9:     Continue to next iteration
10:  end if
11: end if
12:  $\theta^{(i)} \leftarrow \theta^{(i-1)}$ 
13: end for
```

Algorithm 4 SMC-ABC algorithm producing N samples from the approximate posterior distribution $p_{d,\epsilon}(\theta|y_0)$. Here $q_t(\theta|\theta_{t-1}^{(i)}) = N(\theta|\theta_{t-1}^{(i)}, \Sigma_{t-1})$.

Require: Specify a decreasing sequence of thresholds $\epsilon_1 \geq \epsilon_2 \geq \dots \geq \epsilon_T$ for T iterations.

```

1: for  $i = 1$  to  $N$  do
2:   repeat
3:     Generate  $\theta$  from the prior  $p(\cdot)$ 
4:     Generate  $y_\theta$  from the simulator
5:   until  $d(y_\theta, y_0) \leq \epsilon_1$ 
6:    $\theta_1^{(i)} \leftarrow \theta$ 
7:    $\omega_1^{(i)} \leftarrow 1/N$ 
8: end for
9:  $\Sigma_1 \leftarrow 2 \text{Cov}(\theta_1)$  {Twice the empirical variance}
10:
11: for  $t = 2$  to  $T$  do
12:   for  $i = 1$  to  $N$  do
13:     repeat
14:       Draw  $\theta^*$  from among  $\theta_{t-1}$  with probabilities  $\omega_{t-1}$ 
15:       Generate  $\theta$  from  $\mathcal{N}(\theta^*, \Sigma_{t-1})$ 
16:       Generate  $y_\theta$  from the simulator
17:     until  $d(y_\theta, y_0) \leq \epsilon_t$ 
18:      $\theta_t^{(i)} \leftarrow \theta$ 
19:      $\omega_t^{(i)} \leftarrow p(\theta)/(\sum_{k=1}^N \omega_{t-1}^{(k)} \mathcal{N}(\theta|\theta_{t-1}^{(k)}, \Sigma_{t-1}))$  {weights can be scaled with a constant}
20:   end for
21:    $\Sigma_t \leftarrow 2 \text{Cov}(\theta_t)$  {Twice the empirical variance}
22: end for

```
